

類似特許分析にみる、特許分析へのAI活用の現状と可能性

AI for Prior Art Search, Current Problems and Possibilities



日本アイ・ビー・エム株式会社 東京基礎研究所リサーチ・スタッフ・メンバー

鈴木 祥子

2004年に日本アイ・ビー・エム株式会社に入社。東京基礎研究所で数理解析のチームを経て、現在はテキスト分析のチームに所属し、特許文書や技術文書、ビジネス文書などの分析に従事している。博士(理学)。

✉ e30126@jp.ibm.com

1 はじめに

近年のAIの発展は目覚ましく、種々の活用への期待値が高まっている。これは特許関連業務においても同様である。AI適用が進む中、特許分析の専門家の業務はどうなっていくだろうか。本稿では、AIを広く機械的処理と捉えた上で、特に類似特許検索を例にとり、特許分析におけるAI活用の可能性について論じる。

特許関連業務におけるAI活用の期待値の高さにはいくつかの理由がある。第一に、大量の特許データを目的に合わせて正確に分析することに対する高いニーズが理由として挙げられる。特許出願数は世界的に年々増加しており、このような状況において特許データの活用がビジネス上益々重要になってきている。第二の理由は、特許分析の難度の高さである。分析には特許情報自体の専門性および各技術分野の専門性が要求される。また、特許が法的効力を持つという性質上、先行文献調査や無効資料調査などの分析は、その重要性の高さから調査や判断に時間がかかることが多い。AIにより、従来人手で調査しきれなかった大量データを正確に分析することで専門家の分析がより早くより包括的になるのであれば、大きな価値である。特許分析のAI活用への期待値が高い第三の理由は、特許データの使いやすさである。一般にAIのトレーニングには大量のデータを必要とする。実際にAIをなんらかの業務に適用しようと考えたとき、

システム上の都合や著作権による制限、あるいはそもそも利用可能な形式で保存されていない、といった事情により、分析したい十分な量のデータが得られないケースは非常に多い。これに対し、特許データは特許庁の主導の下、(特に日本では)早い時期からテキストデータが電子化されている。またメタデータである書誌情報も、出願日や公開日などのタイムスタンプ、出願人や発明者情報、IPCやFタームといった分野コードなど多種に渡り整備されている。更には審査情報や審判情報など、専門家の判断が含まれる各種データも存在し、学習データとしてAIで利用できる可能性が高い。尚且つ重要なのは、これらの特許情報自体が公共のデータであり入手・分析・活用に対する著作権的な制限がないことである。

実際に翻訳や分類などトレーニングに利用できるデータが大量にある分析業務を中心に特許分野へのAI活用例は広がりを見せている。今後も様々な切り口から活用例は増えていくことが予想される。しかしながら(これはどの分野に適用する際にもいえることだが)、特許という特殊なデータにおいて、他分野で活用されているAIの手法がそのまま適用できるとは限らない。2章では数ある特許関連業務のうち特に類似特許検索においていくつかの例を挙げて、その特殊性を説明する。3章では特許の特殊性に対して、どのようなアプローチが有効か、技術面やユーザビリティの点から議論する。



2 類似特許検索における課題

様々な特許分析業務の中でも、重要な位置を占めるのが類似特許検索である。本章では特許分析業務にAIを活用する際の適用先として、類似特許検索を例に、特許の特殊性とそこから生じる困難を述べる。

類似特許検索は、出願前検索、出願時の先行技術調査、パテントクリアランス、無効資料調査、など多様な目的に応じて様々な局面で行われる。目的によって、ユーザーが求める検索結果も精度も異なるため、本来であれば類似特許検索は目的に応じて異なる仕組みを持つのが望ましいと言える。本章以降では主に先行技術調査を中心に話を進める。それ以外の類似特許検索は、分野や企業ごとのノウハウが存在し、また正解データの入手が部外者には困難なためである。

2.1 一般の文書検索における課題とアプローチ

特許文書に限らず一般の文書検索には様々な課題があり、多くのアプローチが提案されている。本節では、類似文書検索一般における課題を整理し、従来手法や近年多く用いられるアプローチについて簡単に述べる。特に以下の課題について個別に議論する。

- a. クエリとなる文書の重要箇所の特定
- b. クエリ文書とマッチしたい正解文書との間の表現の差異の解消
- c. (評価用、機械学習のための) 正解データ

aの重要箇所の特定については、従来よりキーワードごとの重みが各種提案されておりチューニングされている^[1]が、文書中のどの箇所が重要かは、文書の種類や何を類似とみなすかといった状況ごとに異なる。文書が複数部分に分割されている場合、明示的に部分ごとに重みをつけることも可能である。また、正解データが存在する場合にはランキング学習という手法^[2]が適用可能である。ランキング学習とは、各クエリ文書に対して類似度(適合度)のランキングのついた文書集合の組が大量に存在したとき、これを教師データとして、クエリ文書と検索対象文書から特徴量を算出し、特徴量ごとの重みを学習することで、教師データのランキングを再現するようなモデルを構築する手法である。このランキング学習を利用して文書の各部分の重みを学習するアプロー

チも考えられる。例えば、ある文書集合の検索において、文書の詳細な内容よりもタイトルが一致する方が類似と判断されることが多い場合には、タイトル中のキーワードの重みが大きくなるように学習が行われる。

bの表現の差異の解消は非常に難しい問題である。これは、例えば“今朝は早起きした”と“今日は早朝に目覚めた”という2つの文(それぞれ語集合{“今朝”, “早起きする”}, {“今日”, “早朝”, “目覚める”})を持つ)が、語としては共通していないにもかかわらず人間にとっては類似の意味を持つという問題である。検索においてはクエリとなる入力文書に対して、意味的に類似している文書のランクが上位にあることが望まれるが、従来の、語が一致するかどうかをみる手法では限界があった。クエリ拡張とは、検索におけるクエリを既存の類義語辞書やユーザーフィードバックなどの方法で類義語を含めたキーワード集合に拡張する手法であり、上記検索における表現の差異を解消するために利用されてきた。近年では、これらのアプローチに加え、word2vec^[3]に代表されるword embeddingを用いた語のベクトル表現を用いたアプローチが盛んに適用されている。これらの手法は、各語をある空間上の1点として表現している。この空間が意味の情報を持っているため、語同士の意味的近さが数値的に表現できるようになった。このため、表層上の表現が異なっても類似の文書が見つかる期待が高まってきている。語のベクトル表現だけでなく、文書のベクトル表現手法も提案され、例えば文書のベクトル表現を出力するdoc2vec^[4]では直接文書ベクトルから類似度を算出できる。また短い文書同士の類似度評価のタスクにおいては、従来手法のキーワードごとの重みとword embeddingによる類似度を組み合わせたアプローチも提案されている^[5,6]。

cの正解データについて、上述したように、クエリ文書とランキングのついた検索対象の文書集合の組が多く存在していた場合にはランキング学習^[2]を行うことが出来る。当然ながら、学習の効果が出るためにはこのような正解データは正解ラベルの付け間違いなどのノイズの少ないものである必要がある。また、ランキング学習の適用有無に関わらず、検索精度を評価する上では正解データが常に必要である。しかし現実にはこのような正解データを入手するにはコストがかかる。多くの検索エンジンではクリックログなどを利用して擬似的にランキ

ングを取得している。

2.2 類似特許検索固有の課題

類似文書検索に対し、前節でみたような様々なアプローチがあるものの、類似特許検索には特許固有の困難が存在する。本節では、特許固有の困難を、特許文書の記述の特異性に起因するものと、類似特許検索の正解データの扱いに起因するものとに分けて説明する。

2.2.1 特許文書の記述の特殊性

特許文書は二重の意味で専門的である。1つは法的文書という特性からくる専門性、もうひとつは発明の属する技術分野の専門性である。このため、特許文書は非専門家にとって非常に可読性の低いものとなっている。

特許文書のうち、特に請求項は法的効力を持つため、記述方式が定められている。また法的文書独特の言い回しが存在したり、記述の曖昧性がないよう複雑な内容を厳密に記述するため、文体が一般のテキストと大きく異なり特殊である。例えば特許の請求項には、独立請求項および従属請求項が存在し、請求項間に依存関係がある。

1つの請求項内は複数構成要素に分解され、構成要素間の依存関係が構造として存在する。また、並列構造も多用され、時に入れ子構造になることもある。このような特許の記述の特殊性から、2.1節で述べた通常のテキストで行われるキーワード重み付けが必ずしも当てはまらないことが容易に想像できる。また同じく2.1節で述べた embedding の方法は、語や文の表層上の表現が異なっても類似度を算出できる優れた手法だが、一般的な embedding の手法においては、テキスト内の一定 window 幅内の語の共起を利用している。従って特許請求項のように複雑な構造をもつテキストについて同様のアプローチを行っても、元のテキストの持つ精緻な意味を汲み取ることは難しい。

また、技術的な意味で専門的という点も特許文書の分析において困難となりうる。技術的専門用語は、そもそも一般に概念が複雑で難解である。また、数百万件の特許文書があったとしても、ある分野におけるある専門用語を含む文書数は必ずしも多いとはいえない。例えば、“ホウ化バナジウム”を含む国内文献、あるいは“連体形”を含む国内文献は J-PlatPat^[7] の全文対象の検索ではどちらも数百件程度のヒットであり文献の母数を考える

とかなり少ない数である。一般に、ニューラルネットを含む広く統計的な処理をするにあたり、出現頻度が低い語に対するなんらかの推定の確からしさは低くなる。一方で特許文書において、利用されている専門用語は重要な意味を持つことが多い。このため特許文書の意味を AI が正確に推定することは非常にチャレンジングといえる。専門用語を含む文書の分析には、辞書の利用が非常に有用である。しかし、特許に現れる多様な専門語をカバーする辞書を用意するのは難しい。専門家による整備が必要となるため、一部の分野を除いては体系的な辞書は存在しない。

特許分析の実際の業務においても、これら特許文書の特殊性から調査や審査に非常に時間がかかる。これは何を意味するかというと、機械学習に必要な人手による注釈（アノテーション）付与のコストが高い、ということである。多数の非専門家が大量の画像に対して付与したアノテーションを学習することでニューラルな画像解析が進展したのとは、大きく異なる状況である。

2.2.2 類似特許検索の正解データの特殊性

ここで特許庁の審査における出願時の先行技術調査の特殊性について考える。審査における先行技術調査は、各社ごとに行われる出願前調査やパテントクリアランス、無効資料調査と異なり、専門家の作成した正解データが多数入手できる、という点が異なる。ここで、正解データとは、クエリ文書である各本願請求項のうち、審査の過程で一旦新規性・進歩性なしとして拒絶されたものに対して、拒絶理由として引用された引用文献、と定義する。これらの情報は審査の経過情報から容易に入手できるため、検索精度の評価データとして利用できるほか、ランキング学習時の教師データとしての利用も期待できる。

審査における先行技術調査が、通常の種類文書検索と異なるのは、審査においては類似と判断されるための明確なルールが存在する点である。拒絶理由となる引用文献は本願のある請求項の各構成要素を含んだ文献である必要がある。本願請求項に含まれる全構成要素が1つの文献、あるいは複数の文献の組み合わせで実現できると判断されたとき、本願は拒絶される。このため、類似とされた文献は文書全体が本願と似ている必然性はない。どこか一箇所でも本願の構成要素を表す記述があれば、

それが類似の根拠となるのである。引用文献のどの箇所によって拒絶されたか、についての情報は拒絶理由通知書の中身を読む必要がある。拒絶理由の引用文献情報は一括で大量に入手することが可能である。しかし拒絶理由通知書に書かれた記述内容自体を大量に取得するのは現時点で非常に困難である。

また、上記拒絶理由の引用文献を、通常の類似文書検索の正解データと同等に扱おうとしたときには注意が必要である。これは先行技術調査においては、クエリとなる本願請求項を拒絶するための文献が1つでもあればよい、というルールがあるためである。ある本願を拒絶するために検索対象の文献集合中のすべての類似文献を並べる必要はないのである。そのため拒絶理由となる引用文献は、もしかしたら数多く存在したかもしれない類似文献のうちの一部でしかない。上記の理由から、特許庁の審査の経過情報から得られる引用文献情報だけでは、類似文書検索における完全な正解データとは言えないことが分かる。ランキング学習での教師データとして引用文献を利用する際、あるいは検索精度の評価に引用文献を利用する際には、この点に留意する必要がある。

一方で、検索対象のすべての文献に対し、完全な類似度のランキングを付与するのは非常にコストが高い。また、各検索対象の文献のどの箇所が本願請求項のどの構成要素に対応するかを判断するには専門的な知識が必要であり、人手で精度の高い正解データを大量に用意するには限界がある。大規模な検索精度の評価を行う際や、大量正解データを用いたランキング学習を行う際には、拒絶理由の引用文献の性質を踏まえて、どのような正解データを利用するか考える必要がある。

3 類似特許検索へのアプローチ

前章でみたように、類似特許検索には特許特有の困難がある。本章ではそのいくつかについて、どのようなアプローチが考えられるか議論する。

3.1 特許文書構造解析

特許は多くの構造を持つ文書である。明細書においては、文章が背景技術、発明を解決しようとする課題、具体例などの大きな項目に分割して記述されていることが多い。前章までに挙げたように、法的効力を持つ特許請

求項は更に強固な構造を持っており、請求項間の従属関係や、各請求項の構成要素、更には構成要素依存関係が存在する。また複雑な係り受け関係や入れ子状の並列構造が存在することも多い。これらの構造には意味があり、構造を正しく抽出できることが特許の内容理解そのものに役立つと言える。

特許文書の複雑な構造を分析する試みはいくつか存在する。請求項内の構成要素の分解については、手がかり語や句読点、POS タグなどを用いる手法が知られている^[8,9,10]。構成要素の分解は、内容理解に役立つ他、類似特許検索で構成要素単位のキーワード抽出などに利用できると考えられる。また、構成要素の分解、請求項間依存関係、及び構成要素の依存関係の抽出を行うことで、発明の新規性・進歩性に関する箇所の特定が可能となることが示されている^[11]。ここで例として特開 2017-219937 の請求項の構造の一部を図1に示す。実線の青枠で囲まれているのが請求項1の各構成要素、点線の青枠で囲まれているのが請求項1に従属する請求項、赤線で囲まれている部分が請求項1で新規性・進歩性に関する箇所と特定された要素である。構成要素間の依存関係、および従属請求項との依存関係は矢印で示されている。この手法は、特許請求項の戦略的な記述方法を前提とした手法である。一般に、従属請求項は上位の請求項の請求範囲を限定することで審査や審判後もある程度の権利の範囲を生き残らせようとして記述される。このため、従属請求項によって限定される箇所は発明の最も重要な部分である可能性が高い。多くの場合、これは発明の新規性・進歩性に関する箇所と言える。また、請求項内構成要素にも依存関係が存在するが、これは前述の構成要素を後続の構成要素が限定するために生じるケースが多々ある。限定は権利範囲を狭めるものであるから、発明に不必要な限定は通常行われない。このことから、各請求項中の新規性・進歩性に関する部分については更なる限定は避ける傾向にあるといえる。このような2つの仮定に基づき、請求項の構造解析を利用して新規性・進歩性に関する箇所の推定を行うことが可能になる。実際に既存の重要キーワード抽出手法と比較して、この提案手法が新規性・進歩性に関する部分を精度よく抽出することが確かめられている。

これらの構成要素分解、あるいは新規性・進歩性に関する箇所の抽出は類似特許検索において本願請求項の重

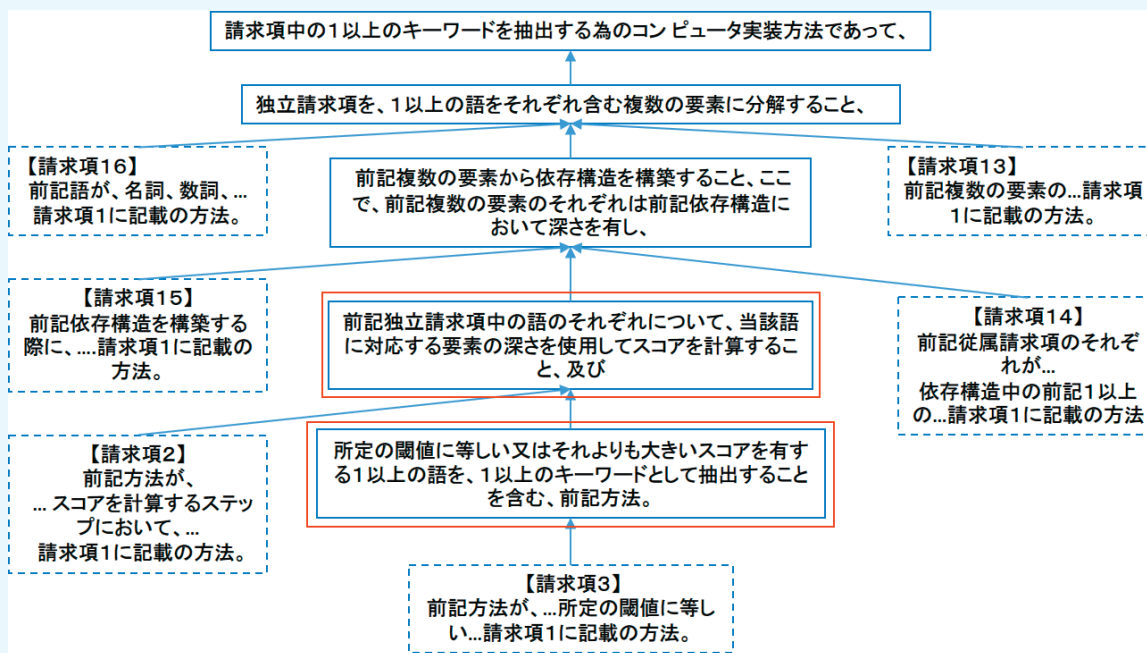


図1 請求項の構造の例

要箇所の特定に役立つと考えられる。また、類似特許検索に限らず、依存関係や構成要素の分解、及び新規性に関する語の抽出は、請求項の可読性を大幅に向上させる。

一方、技術用語の専門性から来る困難についてはどのようなアプローチが考えられるだろうか。出現頻度が低い語に対して出来る限りの情報を取得したいときに有効と考えられるのは、文書中の明示的な説明の利用である^[12]。各特許文書において、請求項の各構成要素についてのより具体的な内容が明細書の詳細説明に記載されている。これらの具体的な説明との紐付けにより、類義語や下位概念候補を見つけることが可能である。また、手がかり語を用いた明示的な関係抽出も、出現頻度は低いながらも正確に語の関係を取得できる。例えば、“通信インターフェース120は、例えば、イーサネット・カードであり、”という記述があったとき、この文書においては“通信インターフェース”の下位概念として“イーサネット・カード”を想定しているということが分かる。このようにして得られた語の関係と、既存辞書、更にはembeddingで得られる関係とを組み合わせることで、語の意味をより正確に捉えることが可能になる。

抽出された特許の構造を利用して専門用語の情報を得ることも考えられる。ここでは、先に述べた請求項間の依存関係から、上位概念-下位概念の候補を抽出することを考えてみる。例えば上位請求項に対して熱可塑性樹脂に関する記述が記載され、下位請求項に対して“熱可

塑性樹脂が、ポリエステル系樹脂、ポリスチレン系樹脂、フッ素系樹脂およびポリメチルペンテン樹脂からなる群から選ばれる”という記述があったとする。このような構造から、上位概念に“熱可塑性樹脂”、下位概念に“ポリエステル系樹脂”、“ポリスチレン系樹脂”、“フッ素系樹脂”、“ポリメチルペンテン樹脂”などの語が候補として挙げられ、上位の“熱可塑性樹脂”と組み合わせが相対的に多い場合には、“熱可塑性樹脂”と下位概念候補との間に上位-下位の関係が存在すると推測できる。図2に、請求項の構造から、上位・下位の関係候補を抽出する例(特開2015-088333)を挙げる。ここでは、請求項1の“凹凸”に対し、請求項3で“高低差”が“0.5~4.0 μm”であることが分かる。同様に請求項1の“炭素材料”の具体例が、請求項4で“カーボンブラック”、“黒鉛”、“カーボンナノチューブ”であると記述されていることが分かる。このような関係を大量に集め、統計的に確からしい関係を選択することで、特許文書からの上位概念・下位概念の候補の抽出が可能になる。

また並列構造を抽出することで知見を得ることも可能である。例えば“機能性基板と、剥離層と、接着剤層と、キャリア基板と、を順に含む仮接合基板システムを形成する方法であって、”という記述からは、“機能性基板”、“剥離層”、“接着剤層”、“キャリア基板”が単純に共起するだけでなく順番も重要であることが分かる。このような知見が複数文書から抽出されれば、頻度は低くても情

【請求項1】集電体上に設けられたカーボンコート層であって、炭素材料による凹凸が活材と接触する側に設けられていることを特徴とする、カーボンコート層。

【請求項2】カーボンコート層を作成するために用いる塗料であって、炭素材料を含み、集電体上への塗工によって前記炭素材料による凹凸が活材と接触する側に形成されることを特徴とする、塗料。

【請求項3】凹凸の高低差は0.5~4.0 μ mである、請求項1に記載のカーボンコート層または請求項2に記載の塗料。

【請求項4】炭素材料は、カーボンブラック、黒鉛およびカーボンナノチューブからなる群より選ばれた少なくとも一種である、請求項1または3に記載のカーボンコート層または請求項2または3に記載の塗料。

【請求項5】異なる粒度の炭素材料が組み合わせられている、請求項1、3および4のいずれか一項に記載のカーボンコート層または請求項2~4のいずれか一項に記載の塗料。

【請求項6】炭素材料は、平均粒径300~700nmの粒度のカーボンブラックをA、平均粒径800~1500nmの粒度のカーボンブラックをBとしたとき、A:B=100:10~300の重量比率で組み合わせられている、請求項5に記載のカーボンコート層または塗料。

【請求項7】請求項1および3~6のいずれか一項に記載のカーボンコート層を有する、集電体または前記集電体を含む電池。

【請求項8】炭素材料による凹凸が活材と接触する側に設けられているカーボンコート層を集電体上に形成する方法であって、以下の工程：炭素材料が含まれる塗料を用意すること、および前記塗料を集電体上に塗布することを含むことを特徴とする、方法。

【請求項9】塗料が請求項2~6のいずれか一項に記載の塗料である、請求項8に記載の方法。

図2 請求項構造を利用した上位・下位概念候補の抽出例

報の精度は上がる。

3.2 類似特許検索へのアプローチ

前節では特許文書の詳細な記述、あるいは複雑な構造を利用して、語の意味や文書の意味を捉える試みを見てきた。このようにして抽出した各種情報は、ハイライトや表示の工夫により特許文書の可読性の向上に役立つ他、類似特許検索へも積極的に利用できる。例えば、本願請求項をクエリとした類似特許検索を行う際には、構成要素分解、あるいは新規性・進歩性に関する箇所の抽出などが2.1節のaで挙げた本願請求項の重要箇所の特定に役立つと考えられる。具体的には、重要箇所の明示的な重み付けや、ランキング学習時の特徴量としてデザインし学習により重み付けを行う、といった利用をすることで、構造のないフラットなテキストとして検索するよりも本願の意味を捉えた検索が可能になる。

また、語の意味の推定精度が上がることで、2.1節のbで挙げたような表層上の語の表現が異なる場合には、既存辞書のみを用いた検索よりも精度が高くなることが期待される。2章でみたようにこのような場合にもembeddingの類似度を用いて文章同士の類似度を計算する手法が提案されている^[5,6]。しかし緻密な比較を行

う場合には、ユーザーに類義語や上位概念、下位概念の候補語を提示し、フィードバックを用いてクエリ拡張を行う方が、ユーザーの意図に沿った効率的な検索を実現できると筆者は考える。

多様な検索目的に対応するための類似特許検索システムとして、ユーザーにとって使いやすく、かつ特許文書の複雑な内容を把握した精度の良い検索を実現するためには、ユーザー側のフィードバックとAIによるサポートとを組み合わせた仕組みが重要である。例えばユーザー向けの機能としては以下が考えられる。

- インputの本願請求項に対して、重要箇所や重要キーワード候補を表示
- ノイズの少ないクエリ拡張の候補語の提示
- ユーザーによるフィードバック（重要キーワード選択、重要箇所重み付け、クエリ拡張時の選択語）を考慮した検索

これらを実現するために、3.1節で解説した請求項構造分析、語のembedding、辞書、手がかり語による関係抽出、構造による関係抽出などのAI機能が背後に必要となる。また、フィードバックを含む最終的なインプットに対し、ランキング学習を行い検索精度を上げる工夫も必要となる。この時の学習データとしては、最初は拒

絶理由の引用文献情報が考えられるが、検索システム上に検索結果のフィードバック機能を付与することで、最終的にはユーザー側の評価自体を学習データとすることが望ましい。このように検索結果のフィードバックを生かすことで 2.1 節の c で挙げた正解データの精度向上自体が期待できる。また、ユーザーのクエリ拡張の選択ログが溜まってくれば、よりよい候補語抽出を学習することも可能である。もちろん、今後拒絶理由通知書の中身が大量に一括で取得できるようになるなど、特許データ周辺の環境自体が変化してきた場合には最大限活用したい。

検索結果の表示に対しても種々の工夫が考えられる。類似特許検索において、正解とされるデータにはいくつかのパターンがある可能性がある。例えば、本願の技術が極めて当たり前のものであった場合、拒絶理由の引用文献は、数ある類似文書の中から、有名でよく引用される文献が選択されることが多い。また、本願と同一発明者による別文献で拒絶されるケースも多くある。これらの複数パターンに対応するためには、文献の類似性の定義を変えた方がよい可能性がある。そのため、パターンに応じた学習データを複数用意し、それぞれのデータから学習されたモデルによる検索結果を複数提供することで、最終的にユーザーにとって望ましい結果が得られる可能性が高まる。

4 おわりに

ここまで、主に類似特許検索に特化して特許分析の特有の課題、および AI によるアプローチの可能性を見てきた。繰り返しになるが、特許文書のような専門性が高い文書の分析は、一般的な AI のアプローチをそのまま適用するよりも、特許文書に特化した処理を併用する方が望ましい。これは純粹に分析精度の観点からも、またユーザビリティの点からも言えることである。現時点での AI 技術による文書理解は発展の途上であり、質のよい学習データが大量に存在するタスクを中心に適用がなされてきた。特許文書における適用も今後一層盛んになると推測されるが、複雑な構造と意味を持つ特許文書の理解には、未だ人間の判断が多く入る必要があると筆者は考える。AI のサポートにより、特許文書の可読性を上げ、また様々な視点での情報を提示することでユーザー

側にとって使いやすいシステムとなり、分析能力自体が向上する。またユーザー側が提示された情報を元にフィードバックを行ったとき、これらのフィードバックを反映する仕組みを作っていくことで、今後活用できる学習データを蓄積していく。このように人間と技術が相互に補い合い進化していくのが今後の特許分析の歩む道ではないだろうか。

* 本稿記載の内容は筆者個人の見解に基づいています。

参考文献

- [1] Stephen Robertson and Hugo Zaragoza, "The Probabilistic Relevance Framework: BM25 and Beyond", Foundations and Trends in Information Retrieval, 2009.
- [2] Tie-Yan Liu. Learning to Rank for Information Retrieval, Springer, 2011.
- [3] Mikolov, T., Chen, K., Corrado, G., and Dean, J. "Efficient estimation of word representations in vector space" , ICLR, 2013.,
- [4] Quoc V Le and Tomas Mikolov. D, "Distributed Representations of Sentences and Documents", ICML, 2014.
- [5] Tom Kenter and Maarten de Rijke, "Short Text Similarity with Word Embeddings", CIKM, 2015
- [6] D. Cer, M. Diab, E. Agirre, I. Lopez-Gazpio, and L. Specia. S, "SemEval-2017 Task 1: Semantic Textual Similarity Multilingual and Cross-lingual Focused Evaluation" , Proc. SemEval, 2017.
- [7] <https://www.j-platpat.inpit.go.jp/>
- [8] Svetlana Sheremetyeva, Sergei Nirenburg, and Irene Nirenburg. "Generating patent claims from interactive input", in Workshop on Natural Language Generation ,1996.
- [9] Peter Parapatics and Michael Dittenbach, "Patent Claim Decomposition for Improved Information Extraction", Workshop on Patent Information Retrieval 2009.
- [10] Akihiro Shinmori, Manabu Okumura, Yuzo Marukawa, and Makoto Iwayama. "Patent Claim Processing for Readability: Structure Analysis and Term Explanation", ACL-2003 Workshop on Patent Corpus Processing, 2003.
- [11] Shoko Suzuki and Hiromichi Takatsuka, "Extraction of Keywords of Novelties from Patent Claims," Coling, 2016
- [12] 奥村 学, "特許情報処理：言語処理のアプローチ", コロナ社, 2012.